

Incrementalizing Random Sketching for Solving Consistent Linear Systems

Vivak Patel

University of Wisconsin – Madison

October 2021

Coauthors & Reference

Vivak Patel, **Mohammad Jahangoshahi**, and **Daniel A. Maldonado**. "An implicit representation and iterative solution of randomly sketched linear systems." *SIAM Journal on Matrix Analysis and Applications* 42.2 (2021): 800-831.

A pre-print is also available on arXiv: 1904.11919.

Outline

1. Overview: Linear System & Randomized Solvers
2. Our Contribution
3. Our Procedure
4. Theory
5. Summary

OVERVIEW

Problem formulation

Our goal is to solve

$$\underbrace{A}_{n \times d} \underbrace{x}_d = \underbrace{b}_n, \quad (1)$$

for which we assume that **at least one solution exists**.

In particular, we are interested in **randomized solvers**, which recast the linear systems problem into a statistical estimation problem and then solve the estimation problem.

Two classes of randomized solvers

1. **Random Sketching Methods.** These methods use a matrix, \mathbf{M} , with (much) fewer rows than \mathbf{A} and solve the problem

$$\min_x \|(\mathbf{M}\mathbf{A})\mathbf{x} - (\mathbf{M}\mathbf{b})\|_2^2, \quad (2)$$

where \mathbf{M} is a specifically structured random matrix, which we call a random sketching matrix.

Two classes of randomized solvers

2. **Base Random Iteration.** These methods use a random vectors, $\{\mathbf{w}_k : k + 1 \in \mathbb{N}\}$ and perform the iteration

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \gamma \mathbf{A}' \mathbf{w}_k [(\mathbf{w}_k' \mathbf{b} - (\mathbf{w}_k' \mathbf{A}) \mathbf{x}_k)], \quad (3)$$

where $\gamma > 0$ is some scalar.

Why randomized solvers?

Random sketching methods have the **lowest computational complexities** for finding an ϵ -accurate solution (in residual) with high probability.

Why randomized solvers?

Random sketching methods have the **lowest computational complexities** for finding an ϵ -accurate solution (in residual) with high probability.

Base random iterations are **cheap-per-iteration** and converge linearly in the number of iterations for appropriately selected $\{W_k\}$.

Why randomized solvers?

Random sketching methods have the **lowest computational complexities** for finding an ϵ -accurate solution (in residual) with high probability.

Base random iterations are **cheap-per-iteration** and converge linearly in the number of iterations for appropriately selected $\{W_k\}$.

Even with high-probability, we would think that the best solution then is the **random sketching** approach as it is the fastest method available.

Challenges with Random Sketching

Challenges with Random Sketching

The number of rows of M depends on constants that are unknown and problem-specific, and is proportional to the reciprocal of ϵ .

Challenges with Random Sketching

The number of rows of \mathbf{M} depends on constants that are unknown and problem-specific, and is proportional to the reciprocal of ϵ .

Therefore, \mathbf{MA} , which can be cheap to compute (if \mathbf{M} is sparse), might still be too expensive to construct and store!

TO SUMMARIZE

Random sketching has the best computational complexity, but we do not know how to choose the size of M and it is nontrivial to store MA .

OUR CONTRIBUTION

Overview & Consequences

We reformulate random sketching to **implicitly** construct MA and **simultaneously** solve the projected system (i.e., $MAx = Mb$).

Overview & Consequences

We reformulate random sketching to **implicitly** construct \mathbf{MA} and **simultaneously** solve the projected system (i.e., $\mathbf{MAx} = \mathbf{Mb}$).

We do **not** need to decide on the size of \mathbf{M} a priori. We can let the size of \mathbf{M} grow implicitly, until, say, some stopping criteria is reached or the system is solved.

Overview & Consequences

We reformulate random sketching to **implicitly** construct MA and **simultaneously** solve the projected system (i.e., $MAx = Mb$).

We do **not** need to decide on the size of M a priori. We can let the size of M grow implicitly, until, say, some stopping criteria is reached or the system is solved.

Additionally, we do **not** need to create and store the matrix MA . We implicitly work with this matrix without constructing it.

BOTTOM LINE

Owing to **our reformulation**, we are able to move towards the **practical** use of random sketching methods to solve actual linear systems.

OUR PROCEDURE

Step 1: Streaming rows of M

Let $w_k \in \mathbb{R}^n$ denote the k^{th} row of a sketching matrix M . Our first requirement is to generate w_k on the fly.

Step 1: Streaming rows of M

Let $\mathbf{w}_k \in \mathbb{R}^n$ denote the k^{th} row of a sketching matrix M . Our first requirement is to generate \mathbf{w}_k on the fly.

Example: Gaussian Sketch. M has independent, identically distributed standard Gaussian entries. Then \mathbf{w}_k is simply an n -dimensional standard Gaussian vector, and each $\{\mathbf{w}_j\}$ are independent.

Step 2: Iterative Solver

Recalling that $\{\mathbf{w}_k\}$ are the rows of our sketching matrix \mathbf{M} , we now work through the iteration

$$\mathbf{x}_{k+1} = \begin{cases} \mathbf{x}_k + \mathbf{S}_k \mathbf{A}' \mathbf{w}_k \frac{\mathbf{w}'_k (\mathbf{b} - \mathbf{A} \mathbf{x}_k)}{\mathbf{w}'_k \mathbf{A} \mathbf{S}_k \mathbf{A}' \mathbf{w}_k} & \mathbf{S}_k \mathbf{A}' \mathbf{w}_k \neq 0 \\ \mathbf{x}_k & \text{otherwise,} \end{cases} \quad (4)$$

and

$$\mathbf{S}_{k+1} = \begin{cases} \mathbf{S}_k - \frac{\mathbf{S}_k \mathbf{A}' \mathbf{w}_k \mathbf{w}'_k \mathbf{A} \mathbf{S}_k}{\mathbf{w}'_k \mathbf{A} \mathbf{S}_k \mathbf{A}' \mathbf{w}_k} & \mathbf{S}_k \mathbf{A}' \mathbf{w}_k \neq 0 \\ \mathbf{S}_k & \text{otherwise,} \end{cases} \quad (5)$$

where $\mathbf{S}_0 = \mathbf{I}_d$ and \mathbf{x}_0 is arbitrary.

Step 2: Iterative Solver

Note, $\{\mathbf{S}_k\}$ are orthogonal projections onto the space perpendicular to the rows of \mathbf{MA} that have already been observed.

Step 2: Iterative Solver

Note, $\{\mathbf{S}_k\}$ are orthogonal projections onto the space perpendicular to the rows of \mathbf{MA} that have already been observed.

In other words, $\mathcal{R}_l := \text{span}[\mathbf{A}'\mathbf{w}_0, \dots, \mathbf{A}'\mathbf{w}_{l-1}]$ then \mathbf{S}_l is an orthogonal projection onto \mathcal{R}_l^\perp .

Step 2: Iterative Solver

Note, $\{\mathbf{S}_k\}$ are orthogonal projections onto the space perpendicular to the rows of \mathbf{MA} that have already been observed.

In other words, $\mathcal{R}_l := \text{span}[\mathbf{A}'\mathbf{w}_0, \dots, \mathbf{A}'\mathbf{w}_{l-1}]$ then \mathbf{S}_l is an orthogonal projection onto \mathcal{R}_l^\perp .

Therefore, as soon as we see the **maximal possible linearly independent rows** of \mathbf{MA} , then we will have solved the system to the highest accuracy possible allowed by \mathbf{M} and \mathbf{A} .

THE CHALLENGE

How do we characterize this maximal set when the rows of M are generated on the fly and they can have an arbitrary (independent, random permutation, adaptive, dependent) structure to previously observed rows of M ?

THEORY

Subspace Characterizations

Let $\mathbf{w} \in \mathbb{R}^n$ be an arbitrary random variable. Define

$$\mathcal{N}(\mathbf{w}) = \text{span} \left[\mathbf{z} \in \mathbb{R}^d : \mathbb{P} [\mathbf{z}'\mathbf{A}'\mathbf{w} = 0] = 1 \right] \quad (6)$$

and

$$\mathcal{R}(\mathbf{w}) = \mathcal{N}(\mathbf{w})^\perp. \quad (7)$$

Subspace Characterizations

Let $\mathbf{w} \in \mathbb{R}^n$ be an arbitrary random variable. Define

$$\mathcal{N}(\mathbf{w}) = \text{span} \left[\mathbf{z} \in \mathbb{R}^d : \mathbb{P} [\mathbf{z}'\mathbf{A}'\mathbf{w} = 0] = 1 \right] \quad (6)$$

and

$$\mathcal{R}(\mathbf{w}) = \mathcal{N}(\mathbf{w})^\perp. \quad (7)$$

Lemma

$\mathcal{R}(\mathbf{w})$ is the smallest subspace of \mathbb{R}^d such that $\mathbb{P} [\mathbf{A}'\mathbf{w} \in \mathcal{R}(\mathbf{w})] = 1$.

Subspace Characterizations

Message: For an arbitrary random variable \mathbf{w}

- $\mathcal{R}(\mathbf{w})$ characterizes the row space of $\mathbf{w}'\mathbf{A}$
- $\mathcal{N}(\mathbf{w})$ characterizes the null space of $\mathbf{w}'\mathbf{A}$
- $\mathcal{V}(\mathbf{w})$ characterizes the deficiency of $\mathbf{w}'\mathbf{A}$ compared to \mathbf{A} .

Iterations to Maximal Set

Now for an arbitrary random variable \mathbf{w} , let $\mathcal{R}(\mathbf{w})$ and $\mathcal{N}(\mathbf{w})$ be defined as before. For (a not necessarily related) set of random variables $\{\mathbf{w}_k\}$, define

$$T = \min\{k \geq 0 : \text{span} [A'\mathbf{w}_0, \dots, A'\mathbf{w}_k] \supset \mathcal{R}(\mathbf{w})\}. \quad (8)$$

Iterations to Maximal Set

Now for an **arbitrary** random variable \mathbf{w} , let $\mathcal{R}(\mathbf{w})$ and $\mathcal{N}(\mathbf{w})$ be defined as before. For (a not necessarily related) set of random variables $\{\mathbf{w}_k\}$, define

$$T = \min\{k \geq 0 : \text{span} [A'\mathbf{w}_0, \dots, A'\mathbf{w}_k] \supset \mathcal{R}(\mathbf{w})\}. \quad (8)$$

Again, we have **not** imposed any relationship between $\mathbf{w}, \mathbf{w}_0, \mathbf{w}_1, \dots$. Therefore, T is quite generally defined. (This is useful when we consider parallel implementations.)

Convergence

Theorem

Let \mathbf{w} be a random variable, and let $\{\mathbf{w}_k\}$ be random variables such that $\mathbb{P}[\mathbf{A}'\mathbf{w}_l \in \mathcal{R}(\mathbf{w})] = 1$ for all $l \geq 0$. Define T as above. On the event $\{T < \infty\}$,

- For any $s \geq T + 1$, $\mathbf{S}_{t+1} = \mathbf{S}_s$ and $\mathbf{x}_{T+1} = \mathbf{x}_s$.
- If $\mathbf{Ax} = \mathbf{b}$ admits a solution \mathbf{x}^* (not necessarily unique), then

$$\mathbf{x}_{T+1} = P_{\mathcal{N}(\mathbf{w})}\mathbf{x}_0 + P_{\mathcal{R}(\mathbf{w})}\mathbf{x}^*. \quad (9)$$

Do we solve the system?

Corollary

Under the settings of the preceding theorem, on the event $\{T < \infty\}$, $\mathbf{Ax}_{T+1} = \mathbf{b}$ if and only if $\mathbf{P}_{\mathcal{V}(w)}\mathbf{x}_0 = \mathbf{P}_{\mathcal{V}(w)}\mathbf{x}^*$.

Do we solve the system?

Corollary

Under the settings of the preceding theorem, on the event $\{T < \infty\}$, $\mathbf{Ax}_{T+1} = \mathbf{b}$ if and only if $\mathbf{P}_{\mathcal{V}(w)}\mathbf{x}_0 = \mathbf{P}_{\mathcal{V}(w)}\mathbf{x}^*$.

- (1) When is $T < \infty$?
- (2) When will $\mathbf{P}_{\mathcal{V}(w)}(\mathbf{x}^* - \mathbf{x}_0) = 0$?

Do we solve the system?

Corollary

Under the settings of the preceding theorem, on the event $\{T < \infty\}$, $\mathbf{Ax}_{T+1} = \mathbf{b}$ if and only if $\mathbf{P}_{\mathcal{V}(w)}\mathbf{x}_0 = \mathbf{P}_{\mathcal{V}(w)}\mathbf{x}^*$.

(1) When is $T < \infty$?

(2) When will $\mathbf{P}_{\mathcal{V}(w)}(\mathbf{x}^* - \mathbf{x}_0) = \mathbf{0}$?

Basically, when is this going to actually work?

When is this going to work?

Both of these questions will depend on how you choose \mathbf{w} , and how you design $\mathbf{w}_0, \mathbf{w}_1, \dots$ for your particular system. This should depend on the **linear system's structure** and **the hardware environment**.

When is this going to work?

Both of these questions will depend on how you choose \mathbf{w} , and how you design $\mathbf{w}_0, \mathbf{w}_1, \dots$ for your particular system. This should depend on the **linear system's structure** and **the hardware environment**.

We have simply stated a **very general theory of convergence** for such methods, and supply specific examples in the paper.

SUMMARY

We restated matrix sketching as a **random orthogonalization procedure** and characterized the convergence for **arbitrary** sampling methodologies. This allows us to **implicitly and incrementally** generate and grow *MA* without storing it explicitly.

THANK YOU

www.vivakpatel.org